

---

# The Importance of Non-Markovianity in Maximum State Entropy Exploration

---

Mirco Mutti<sup>\*1,2</sup> Riccardo De Santi<sup>\*3</sup> Marcello Restelli<sup>1</sup>

## Abstract

In the maximum state entropy exploration framework, an agent interacts with a reward-free environment to learn a policy that maximizes the entropy of the expected state visitations it is inducing. Hazan et al. (2019) noted that the class of Markovian stochastic policies is sufficient for the maximum state entropy objective, and exploiting non-Markovianity is generally considered pointless in this setting. In this paper, we argue that non-Markovianity is instead paramount for maximum state entropy exploration in a finite-sample regime. Especially, we recast the objective to target the expected entropy of the induced state visitations in a single trial. Then, we show that the class of non-Markovian deterministic policies is sufficient for the introduced objective, while Markovian policies suffer non-zero regret in general. However, we prove that the problem of finding an optimal non-Markovian policy is NP-hard. Despite this negative result, we discuss avenues to address the problem in a tractable way and how non-Markovian exploration could benefit the sample efficiency of online reinforcement learning in future works.

## 1. Introduction

Several recent works have addressed *Maximum State Entropy* (MSE) exploration (Hazan et al., 2019; Tarbouriech & Lazaric, 2019; Lee et al., 2019; Mutti & Restelli, 2020; Mutti et al., 2021; Zhang et al., 2021; Guo et al., 2021; Liu & Abbeel, 2021b;a; Seo et al., 2021; Yarats et al., 2021; Mutti et al., 2022; Nedergaard & Cook, 2022) as an objective for unsupervised Reinforcement Learning (RL) (Sutton & Barto, 2018). In this line of work, an agent interacts with a reward-free environment (Jin et al., 2020) in order to maximize an

entropic measure of the state distribution induced by its behavior over the environment, effectively targeting a uniform exploration of the state space. Previous works motivated this MSE objective in two main directions. On the one hand, this learning procedure can be seen as a form of *unsupervised pre-training* of the base model (Laskin et al., 2021), which has been extremely successful in supervised learning (Erhan et al., 2009; 2010; Brown et al., 2020). In this view, a MSE policy can serve as an exploratory initialization to standard learning techniques, such as Q-learning (Watkins & Dayan, 1992) or policy gradient (Peters & Schaal, 2008), and this has been shown to benefit the sample efficiency of a variety of RL tasks that could be specified over the pre-training environment (e.g., Mutti et al., 2021; Liu & Abbeel, 2021b; Laskin et al., 2021). On the other hand, pursuing a MSE objective leads to an even coverage of the state space, which can be instrumental to address the *sparse reward discovery* problem (Tarbouriech et al., 2021). Especially, even when the fine-tuning is slow (Campos et al., 2021), the MSE policy might allow to solve hard-exploration tasks that are out of reach of RL from scratch (Mutti et al., 2021; Liu & Abbeel, 2021b). As we find these premises fascinating, and of general interest to the RL community, we believe it is worth providing a theoretical reconsideration of the MSE problem. Specifically, we aim to study the minimal class of policies that is necessary to optimize a well-posed MSE objective, and the general complexity of the resulting learning problem.

All of the existing works pursuing a MSE objective solely focus on optimizing Markovian exploration strategies, in which each decision is conditioned on the current state of the environment rather than the full history of the visited states. The resulting learning problem is known to be provably efficient in tabular domains (Hazan et al., 2019; Zhang et al., 2020). Moreover, this choice is common in RL, as it is well-known that an optimal deterministic Markovian strategy maximizes the usual cumulative sum of rewards objective (Puterman, 2014). Similarly, Hazan et al. (2019, Lemma 3.3) note that the class of Markovian strategies is *sufficient* for the standard MSE objective. A carefully constructed Markovian strategy is able to induce the same state distribution of any history-based (non-Markovian) one by exploiting randomization. Crucially, this result does not hold only for asymptotic state distributions, but also for state

---

<sup>\*</sup>Equal contribution <sup>1</sup>Politecnico di Milano <sup>2</sup>Università di Bologna <sup>3</sup>ETH Zurich. Correspondence to: Mirco Mutti <mirco.mutti@polimi.it>, Riccardo De Santi <rdesanti@ethz.ch>.

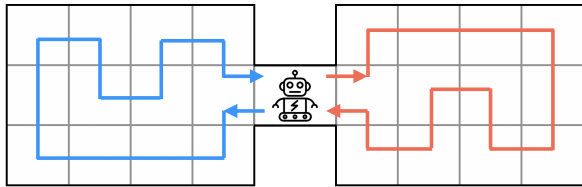


Figure 1. Illustrative two-rooms domain. The agent starts in the middle, colored traces represent optimal strategies to explore the left and the right room.

distributions that are marginalized over a finite horizon (Puterman, 2014). Hence, there is little incentive to consider more complicated strategies as they are not providing any benefit on the value of the entropy objective.

However, the intuition suggests that exploiting the history of the interactions is useful when the agent’s goal is to uniformly explore the environment: If you know what you have visited already, you can take decisions accordingly. To this point, let us consider an illustrative example in which the agent finds itself in the middle of a two-rooms domain (as depicted in Figure 1), having a budget of interactions that is just enough to visit every state within a single episode. It is easy to see that an optimal Markovian strategy for the MSE objective would randomize between going left and right in the initial position, and then would follow the optimal route within a room, finally ending in the initial position again. An episode either results in visiting the left room twice, or the right room twice, or each room once, and all of this outcomes have the same probability. Thus, the agent might explore poorly when considering a single episode, but the exploration is uniform in the average of *infinite trials*. Arguably, this is quite different from how a human being would tackle this problem, i.e., taking intentional decisions in the middle position to visit a room before the other. This strategy leads to uniform exploration of the environment in *any trial*, but it is inherently non-Markovian.

Backed by this intuition, we argue that prior work does not recognize the importance of non-Markovianity in MSE exploration due to an hidden infinite-samples assumption in the objective formulation, which is in sharp contrast with the objective function it is actually optimized by empirical methods, i.e., the state entropy computed over a finite batch of interactions. In this paper, we introduce a new *finite-sample* MSE objective that is akin to the practical formulation, as it targets the expected entropy of the state visitation frequency induced within an episode instead of the entropy of the expected state visitation frequency over infinite samples. In this finite-sample formulation non-Markovian strategies are crucial, and we believe they can benefit a significant range of relevant applications. For example, collecting task-specific samples might be costly in some real-world domains, and

a pre-trained non-Markovian strategy is essential to guarantee quality exploration even in a single-trial setting. In another instance, one might aim to pre-train an exploration strategy for a class of multiple environments instead of a single one. A non-Markovian strategy could exploit the history of interactions to swiftly identify the structure of the environment, then employing the environment-specific optimal strategy thereafter. Unfortunately, learning a non-Markovian strategy is in general much harder than a Markovian one, and we are able to show that it is NP-hard in this setting. Nonetheless, this paper aims to highlight the importance of non-Markovianity to fulfill the promises of maximum state entropy exploration, thereby motivating the development of tractable formulations of the problem as future work.

The contributions are organized as follows. First, in Section 3, we report a known result (Puterman, 2014) to show that the class of Markovian strategies is sufficient for any infinite-samples MSE objective, including the entropy of the induced marginal state distributions in episodic settings. Then, in Section 4, we propose a novel finite-sample MSE objective and a corresponding regret formulation. Especially, we prove that the class of non-Markovian strategies is sufficient for the introduced objective, whereas the optimal Markovian strategy suffers a non-zero regret. However, in Section 5, we show that the problem of finding an optimal non-Markovian strategy for the finite-sample MSE objective is NP-hard in general. Despite the hardness result, we provide a numerical validation of the theory (Section 6), and we comment some potential options to address the problem in a tractable way (Section 7). In Appendix A, we discuss the related work in the MSE literature, while the missing proofs can be found in Appendix B.

## 2. Preliminaries

In the following, we will denote with  $\Delta(\mathcal{X})$  the simplex of a space  $\mathcal{X}$ , with  $[T]$  the set of integers  $\{0, \dots, T - 1\}$ , and with  $v \oplus u$  a concatenation of the vectors  $v, u$ .

**Controlled Markov Process** A Controlled Markov Process (CMP) is a tuple  $\mathcal{M} := (\mathcal{S}, \mathcal{A}, P, \mu)$ , where  $\mathcal{S}$  is a finite state space ( $|\mathcal{S}| = S$ ),  $\mathcal{A}$  is a finite action space ( $|\mathcal{A}| = A$ ),  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is the transition model, such that  $P(s'|a, s)$  denotes the probability of reaching state  $s' \in \mathcal{S}$  when taking action  $a \in \mathcal{A}$  in state  $s \in \mathcal{S}$ , and  $\mu \in \Delta(\mathcal{S})$  is the initial state distribution.

**Policies** A policy  $\pi$  defines the behavior of an agent interacting with an environment modelled by a CMP. It consists of a sequence of decision rules  $\pi := (\pi_t)_{t=0}^{\infty}$ . Each of them is a map between histories  $h := (s_j, a_j)_{j=0}^t \in \mathcal{H}_t$  and actions  $\pi_t : \mathcal{H}_t \rightarrow \Delta(\mathcal{A})$ , such that  $\pi_t(a|h)$  defines

the conditional probability of taking action  $a \in \mathcal{A}$  having experienced the history  $h \in \mathcal{H}_t$ . We denote as  $\mathcal{H}$  the space of the histories of arbitrary length. We denote as  $\Pi$  the set of all the policies, and as  $\Pi^D$  the set of deterministic policies  $\pi = (\pi_t)_{t=1}^\infty$  such that  $\pi_t : \mathcal{H}_t \rightarrow \mathcal{A}$ . We further define:

- **Non-Markovian (NM)** policies  $\Pi_{\text{NM}}$ , where each  $\pi \in \Pi_{\text{NM}}$  collapses to a single time-invariant decision rule  $\pi = (\pi, \pi, \dots)$  such that  $\pi : \mathcal{H} \rightarrow \Delta(\mathcal{A})$ ;
- **Markovian (M)** policies  $\Pi_{\text{M}}$ , where each  $\pi \in \Pi_{\text{M}}$  is defined through a sequence of Markovian decision rules  $\pi = (\pi_t)_{t=0}^\infty$  such that  $\pi_t : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ . A Markovian policy that collapses into a single time-invariant decision rule  $\pi = (\pi, \pi, \dots)$  is called a *stationary* policy.

**State Distributions and Visitation Frequency** A policy  $\pi \in \Pi$  interacting with a CMP induces a  $t$ -step state distribution  $d_t^\pi(s) := \Pr(s_t = s | \pi)$  over  $\mathcal{S}$  (Puterman, 2014). This distribution is described by the temporal relation  $d_t^\pi(s) = \int_{\mathcal{S}} \int_{\mathcal{A}} d_{t-1}^\pi(s', a') P(s | s', a') ds' da'$ , where  $d_t^\pi(\cdot, \cdot) \in \Delta(\mathcal{S} \times \mathcal{A})$  is the  $t$ -step state-action distribution. We call the asymptotic fixed point of this temporal relation the *stationary state distribution*  $d_\infty^\pi(s) := \lim_{t \rightarrow \infty} d_t^\pi(s)$ , and we denote as  $d_\gamma^\pi(s) := (1 - \gamma) \sum_{t=0}^\infty \gamma^t d_t^\pi(s)$  its  $\gamma$ -discounted counterpart, where  $\gamma \in (0, 1)$  is the discount factor. A marginalization of the  $t$ -step state distribution over a finite horizon  $T$ , i.e.,  $d_T^\pi(s) := \frac{1}{T} \sum_{t \in [T]} d_t^\pi(s)$ , is called the *marginal state distribution*. The *state visitation frequency*  $d_h(s) = \frac{1}{T} \sum_{t \in [T]} \mathbb{1}(s_t = s | h)$  is a realization of the marginal state distribution, such that  $\mathbb{E}_{h \sim p_T^\pi} [d_h(s)] = d_T^\pi(s)$ , where the distribution over histories  $p_T^\pi \in \Delta(\mathcal{H}_T)$  is defined as  $p_T^\pi(h) = \mu(s_0) \prod_{t \in [T-1]} \pi(a_t | h_t) P(s_{t+1} | a_t, s_t)$ .

**Markov Decision Process** A CMP  $\mathcal{M}$  paired with a reward function  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is called a Markov Decision Process (MDP) (Puterman, 2014)  $\mathcal{M}^R := \mathcal{M} \cup R$ . We denote with  $R(s, a)$  the expected immediate reward when taking action  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$ , and with  $R(h) = \sum_{t \in [T]} R(s_t, a_t)$  the return over the horizon  $T$ . The performance of a policy  $\pi$  over the MDP  $\mathcal{M}^R$  is defined as the *average return*  $\mathcal{J}_{\mathcal{M}^R}(\pi) = \mathbb{E}_{h \sim p_T^\pi} [R(h)]$ , and  $\pi_{\mathcal{J}}^* \in \arg \max_{\pi \in \Pi} \mathcal{J}_{\mathcal{M}^R}(\pi)$  is called an optimal policy. For any MDP  $\mathcal{M}^R$ , there always exists a deterministic Markovian policy  $\pi \in \Pi_{\text{M}}^D$  that is optimal (Puterman, 2014).

**Extended MDP** The problem of finding an optimal non-Markovian policy with history-length  $T$  in an MDP  $\mathcal{M}^R$ , i.e.,  $\pi_{\text{NM}}^* \in \arg \max_{\pi \in \Pi_{\text{NM}}} \mathcal{J}_{\mathcal{M}^R}(\pi)$ , can be reformulated as the one of finding an optimal Markovian policy  $\pi_{\text{M}}^* \in \arg \max_{\pi \in \Pi_{\text{M}}} \mathcal{J}_{\widetilde{\mathcal{M}}_T^R}(\pi)$  in an extended MDP  $\widetilde{\mathcal{M}}_T^R$ . The extended MDP is defined as  $\widetilde{\mathcal{M}}_T^R := (\widetilde{\mathcal{S}}, \widetilde{\mathcal{A}}, \widetilde{P}, \widetilde{R}, \widetilde{\mu})$ , in which  $\widetilde{\mathcal{S}} \subseteq \mathcal{H}_{[T]} = \mathcal{H}_1 \cup \dots \cup \mathcal{H}_T$ , and  $\widetilde{s} := (\widetilde{s}_0, \dots, \widetilde{s}_{-1})$

corresponds to a history in  $\mathcal{M}^R$  of length  $|\widetilde{s}|$ ,  $\widetilde{\mathcal{A}} = \mathcal{A}$ ,  $\widetilde{P}(\widetilde{s}' | \widetilde{s}, \widetilde{a}) = P(s' = \widetilde{s}'_{-1} | s = \widetilde{s}_{-1}, a = \widetilde{a})$ ,  $\widetilde{R}(\widetilde{s}, \widetilde{a}) = R(s = \widetilde{s}_{-1}, a = \widetilde{a})$ , and  $\widetilde{\mu}(\widetilde{s}) = \mu(s = \widetilde{s})$  for any  $\widetilde{s} \in \widetilde{\mathcal{S}}$  of unit length.

**Partially Observable MDP** A Partially Observable Markov Decision Process (POMDP) (Astrom, 1965; Kaelbling et al., 1998) is described by  $\mathcal{M}_\Omega^R := (\mathcal{S}, \mathcal{A}, P, R, \mu, \Omega, O)$ , where  $\mathcal{S}, \mathcal{A}, P, R, \mu$  are defined as in an MDP,  $\Omega$  is a finite observation space, and  $O : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\Omega)$  is the observation function, such that  $O(o | s', a)$  denotes the conditional probability of the observation  $o \in \Omega$  when selecting action  $a \in \mathcal{A}$  in state  $s \in \mathcal{S}$ . Crucially, while interacting with a POMDP the agent cannot observe the state  $s \in \mathcal{S}$ , but just the observation  $o \in \Omega$ . The performance of a policy  $\pi$  is defined as in an MDP.

### 3. Infinite Samples: Non-Markovianity Does Not Matter

Previous works pursuing maximum state entropy exploration of a CMP consider an objective of the kind

$$\mathcal{E}_\infty(\pi) := H(d^\pi(\cdot)) = - \mathbb{E}_{s \sim d^\pi} [\log d^\pi(s)], \quad (1)$$

where  $d^\pi(\cdot)$  is either a stationary state distribution (Mutti & Restelli, 2020), a discounted state distribution (Hazan et al., 2019; Tarbouriech & Lazaric, 2019), or a marginal state distribution (Lee et al., 2019; Mutti et al., 2021). While it is well-known (Puterman, 2014) that there exists an optimal deterministic policy  $\pi^* \in \Pi_{\text{M}}^D$  for the common average return objective  $\mathcal{J}_{\mathcal{M}^R}$ , it is not pointless to wonder whether the objective in (1) requires a more powerful policy class than  $\Pi_{\text{M}}$ . Hazan et al. (2019, Lemma 3.3) confirm that the set of (randomized) Markovian policies  $\Pi_{\text{M}}$  is indeed sufficient for  $\mathcal{E}_\infty$  defined over asymptotic (stationary or discounted) state distributions. In the following theorem and corollary, we report a common MDP result (Puterman, 2014) to show that  $\Pi_{\text{M}}$  suffices for  $\mathcal{E}_\infty$  defined over (non-asymptotic) marginal state distributions as well.

**Theorem 3.1.** *Let  $x \in \{\infty, \gamma, T\}$ , and let  $\mathcal{D}_{\text{NM}}^x = \{d_x^\pi(\cdot) : \pi \in \Pi_{\text{NM}}\}$ ,  $\mathcal{D}_{\text{M}}^x = \{d_x^\pi(\cdot) : \pi \in \Pi_{\text{M}}\}$  the corresponding sets of state distributions over a CMP. We can prove that:*

- The sets of stationary state distributions are equivalent  $\mathcal{D}_{\text{NM}}^\infty \equiv \mathcal{D}_{\text{M}}^\infty$ ;*
- The sets of discounted state distributions are equivalent  $\mathcal{D}_{\text{NM}}^\gamma \equiv \mathcal{D}_{\text{M}}^\gamma$  for any  $\gamma$ ;*
- The sets of marginal state distributions are equivalent  $\mathcal{D}_{\text{NM}}^T \equiv \mathcal{D}_{\text{M}}^T$  for any  $T$ .*

*Proof Sketch.* For any non-Markovian policy  $\pi \in \Pi_{\text{NM}}$  inducing distributions  $d_t^\pi(\cdot), d_t^\pi(\cdot, \cdot)$  over the states and the

state-action pairs of the CMP, we can build a Markovian policy  $\pi' \in \Pi_M, \pi' = (\pi'_t)_{t=0}^\infty$  through the construction  $\pi'_t(a|s) = d_t^\pi(s, a)/d_t^\pi(s), \forall s \in \mathcal{S}, \forall a \in \mathcal{A}$ . From (Puterman, 2014, Theorem 5.5.1) we know that  $d_t^\pi(s) = d_t^{\pi'}(s)$  holds for any  $t \geq 0$  and  $\forall s \in \mathcal{S}$ . This implies that  $d_\infty^\pi(\cdot) = d_\infty^{\pi'}(\cdot), d_\gamma^\pi(\cdot) = d_\gamma^{\pi'}(\cdot), d_T^\pi(\cdot) = d_T^{\pi'}(\cdot)$ , from which  $\mathcal{D}_{NM}^x \equiv \mathcal{D}_M^x$  follows. See Appendix B.1 for a detailed proof.  $\square$

From the equivalence of the sets of induced distributions, it is straightforward to derive the optimality of Markovian policies for objective (1).

**Corollary 3.2.** *For every CMP, there exists a Markovian policy  $\pi^* \in \Pi_M$  such that  $\pi^* \in \arg \max_{\pi \in \Pi} \mathcal{E}_\infty(\pi)$ .*

As a consequence of Corollary 3.2, there is little incentive to consider non-Markovian policies when optimizing objective (1), since there is no clear advantage to make up for the additional complexity of the policy. This result might be unsurprising when considering asymptotic distributions, as one can expect a carefully constructed Markovian policy to be able to tie the distribution induced by a non-Markovian policy in the limit of the interaction steps. However, it is less evident that a similar property holds for the expectation of final-length interactions alike. Yet, we were able to show that a Markovian policy that properly exploits randomization can always achieve equivalent state distributions w.r.t. non-Markovian counterparts. Note that state distributions are actually *expected* state visitation frequency, and the expectation practically implies an infinite number of realizations. In this paper, we show that this underlying infinite-sample regime is the reason why the benefit of non-Markovianity, albeit backed up by intuition, does not matter. Instead, we propose a relevant finite-sample entropy objective in which non-Markovianity is crucial.

#### 4. Finite Samples: Non-Markovianity Matters

In this section, we reformulate the typical maximum state entropy exploration objective of a CMP (1) to account for a finite-sample regime. Crucially, we consider the expected entropy of the state visitation frequency rather than the entropy of the expected state visitation frequency, which results in

$$\mathcal{E}(\pi) := \mathbb{E}_{h \sim p_T^\pi} [H(d_h(\cdot))] = - \mathbb{E}_{h \sim p_T^\pi} \mathbb{E}_{s \sim d_h} [\log d_h(s)]. \quad (2)$$

We note that  $\mathcal{E}(\pi) \leq \mathcal{E}_\infty(\pi)$  for any  $\pi \in \Pi$ , which is trivial by the concavity of the entropy function and the Jensen's inequality. Whereas (2) is ultimately an expectation as it is (1), the entropy is not computed over the infinite-sample state distribution  $d_T^\pi(\cdot)$  but its finite-sample realization  $d_h(\cdot)$ . Thus, to maximize  $\mathcal{E}(\pi)$  we have to find a policy inducing high-entropy state visits within a single trajectory rather

than high-entropy state visits over infinitely many trajectories. Crucially, while Markovian policies are as powerful as any other policy class in terms of induced state distributions (Theorem 3.1), this is no longer true when looking at induced trajectory distributions  $p_T^\pi$ . Indeed, we will show that non-Markovianity provides a superior policy class for objective (2). First, we define a performance measure to formally assess this benefit, which we call the *regret-to-go*.<sup>1</sup>

**Definition 4.1** (Expected Regret-to-go). *Consider a policy  $\pi \in \Pi$  interacting with a CMP over  $T - t$  steps starting from the trajectory  $h_t$ . We define the expected regret-to-go  $\mathcal{R}_{T-t}$ , i.e., from step  $t$  onwards, as*

$$\mathcal{R}_{T-t}(\pi, h_t) = H^* - \mathbb{E}_{h_{T-t} \sim p_{T-t}^\pi} [H(d_{h_t \oplus h_{T-t}}(\cdot))],$$

where  $H^* = \max_{\pi^* \in \Pi} \mathbb{E}_{h_{T-t}^* \sim p_{T-t}^{\pi^*}} [H(d_{h_t \oplus h_{T-t}^*}(\cdot))]$  is the expected entropy achieved by an optimal policy  $\pi^*$ . The term  $\mathcal{R}_T(\pi)$  denotes the expected regret-to-go of a  $T$ -step trajectory  $h_T$  starting from  $s \sim \mu$ .

The intuition behind the regret-to-go is quite simple. Suppose to have drawn a trajectory  $h_t$  upon step  $t$ . If we take the subsequent action with the (possibly sub-optimal) policy  $\pi$ , by how much would we decrease (in expectation) the entropy of the state visits  $H(d_{h_t}(\cdot))$  w.r.t. an optimal policy  $\pi^*$ ? In particular, we would like to know how limiting the policy  $\pi$  to a specific policy class would affect the expected regret-to-go and the value of  $\mathcal{E}(\pi)$  we could achieve. The following theorem and subsequent corollary, which constitute the main contribution of this paper, state that an optimal non-Markovian policy suffers zero expected regret-to-go in any case, whereas an optimal Markovian policy suffers non-zero expected regret-to-go in general.

**Theorem 4.2** (Non-Markovian Optimality). *For every CMP  $\mathcal{M}$  and trajectory  $h_t \in \mathcal{H}_{[T]}$ , there exists a deterministic non-Markovian policy  $\pi_{NM} \in \Pi_{NM}^D$  that suffers zero regret-to-go  $\mathcal{R}_{T-t}(\pi_{NM}, h_t) = 0$ , whereas for any  $\pi_M \in \Pi_M$  we have  $\mathcal{R}_{T-t}(\pi_M, h_t) \geq 0$ .*

**Corollary 4.3** (Sufficient Condition). *For every CMP  $\mathcal{M}$  and trajectory  $h_t \in \mathcal{H}_{[T]}$  for which any optimal Markovian policy  $\pi_M \in \Pi_M$  is randomized (i.e., stochastic) in  $s_t$ , we have strictly positive regret-to-go  $\mathcal{R}_{T-t}(\pi_M, h_t) > 0$ .*

The result of Theorem 4.2 highlights the importance of non-Markovianity for optimizing the finite-sample MSE objective (2), as the class of Markovian policies is dominated by the class of non-Markovian policies. Most importantly, Corollary 4.3 shows that non-Markovian policies are strictly

<sup>1</sup>Note that the entropy function does not enjoy additivity, thus we cannot adopt the usual expected cumulative regret formulation in this setting.



better than Markovian policies in any CMP of practical interest, i.e., those in which any optimal Markovian policy has to be randomized (Hazan et al., 2019) in order to maximize (2). The intuition behind this result is that a Markovian policy would randomize to make up for the uncertainty over the history, whereas a non-Markovian policy does not suffer from this partial observability, and it can deterministically select an optimal action. Clearly, this partial observability is harmless when dealing with the standard RL objective, in which the reward is fully Markovian and does not depend on the history, but it is instead relevant in the peculiar MSE setting, in which the objective is a concave function of the state visitation frequency. In the following section, we report a sketch of the derivation underlying Theorem 4.2 and Corollary 4.3, while we refer to the Appendix B.2 for complete proofs.

#### 4.1. Regret Analysis

To the purpose of the regret analysis, we will consider the following assumption to ease the notation.<sup>2</sup>

**Assumption 1** (Unique Optimal Action). *For every CMP  $\mathcal{M}$  and trajectory  $h_t \in \mathcal{H}_{[T]}$ , there exists a unique optimal action  $a^* \in \mathcal{A}$  w.r.t. the objective (2).*

First, we show that the class of deterministic non-Markovian policies is sufficient for the minimization of the regret-to-go, and thus for the maximization of (2).

**Lemma 4.4.** *For every CMP  $\mathcal{M}$  and trajectory  $h_t \in \mathcal{H}_{[T]}$ , there exists a deterministic non-Markovian policy  $\pi_{\text{NM}} \in \Pi_{\text{NM}}^{\text{D}}$  such that  $\pi_{\text{NM}} \in \arg \max_{\pi \in \Pi_{\text{NM}}} \mathcal{E}(\pi)$ , which suffers zero regret-to-go  $\mathcal{R}_{T-t}(\pi_{\text{NM}}, h_t) = 0$ .*

*Proof.* The result  $\mathcal{R}_{T-t}(\pi_{\text{NM}}, h_t) = 0$  is straightforward by noting that the set of non-Markovian policies  $\Pi_{\text{NM}}$  with arbitrary history-length is as powerful as the general set of policies  $\Pi$ . To show that there exists a deterministic  $\pi_{\text{NM}}$ , we consider the extended MDP  $\tilde{\mathcal{M}}_T^R$  obtained from the CMP  $\mathcal{M}$  as in Section 2, in which the extended reward function is  $\tilde{R}(\tilde{s}, \tilde{a}) = H(d_{\tilde{s}}(\cdot))$  for every  $\tilde{a} \in \tilde{\mathcal{A}}$ , and every  $\tilde{s} \in \tilde{\mathcal{S}}$  such that  $|\tilde{s}| = T$ , and  $\tilde{R}(\tilde{s}, \tilde{a}) = 0$  otherwise. Since a Markovian policy  $\tilde{\pi}_M \in \Pi_M^{\text{D}}$  on  $\tilde{\mathcal{M}}_T^R$  can be mapped to a non-Markovian policy  $\pi_{\text{NM}} \in \Pi_{\text{NM}}^{\text{D}}$  on  $\mathcal{M}$ , and it is well-known (Puterman, 2014) that for any MDP there exists an optimal deterministic Markovian policy, we have that  $\tilde{\pi}_M \in \arg \max_{\pi \in \Pi_M} \mathcal{J}_{\tilde{\mathcal{M}}_T^R}(\pi)$  implies  $\pi_{\text{NM}} \in \arg \max_{\pi \in \Pi_{\text{NM}}} \mathcal{E}(\pi)$ .  $\square$

<sup>2</sup>Note that this assumption could be easily removed by partitioning the action space in  $h_t$  as  $\mathcal{A}(h_t) = \mathcal{A}_{\text{opt}}(h_t) \cup \mathcal{A}_{\text{sub-opt}}(h_t)$ , such that  $\mathcal{A}_{\text{opt}}(h_t)$  are optimal actions and  $\mathcal{A}_{\text{sub-opt}}(h_t)$  are sub-optimal, and substituting any term  $\pi(a^*|h_t)$  with  $\sum_{a \in \mathcal{A}_{\text{opt}}(h_t)} \pi(a|h_t)$  in the results.

Then, in order to prove that the class of non-Markovian policies is also necessary for regret minimization, it is worth showing that Markovian policies can instead rely on randomization to optimize objective (2).

**Lemma 4.5.** *Let  $\pi_{\text{NM}} \in \Pi_{\text{NM}}^{\text{D}}$  a non-Markovian policy such that  $\pi_{\text{NM}} \in \arg \max_{\pi \in \Pi} \mathcal{E}(\pi)$  on a CMP  $\mathcal{M}$ . For a fixed history  $h_t \in \mathcal{H}_t$  ending in state  $s$ , the variance of the event of an optimal Markovian policy  $\pi_M \in \arg \max_{\pi \in \Pi_M} \mathcal{E}(\pi)$  taking  $a^* = \pi_{\text{NM}}(h_t)$  in  $s$  is given by*

$$\text{Var} [\mathcal{B}(\pi_M(a^*|s, t))] = \text{Var}_{h_s \sim p_t^{\pi_{\text{NM}}}} [\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^*|hs))]],$$

where  $hs \in \mathcal{H}_t$  is any history of length  $t$  such that the final state is  $s$ , i.e.,  $hs := (h_{t-1} \in \mathcal{H}_{t-1}) \oplus s$ , and  $\mathcal{B}(x)$  is a Bernoulli with parameter  $x$ .

*Proof Sketch.* We can prove the result through the Law of Total Variance (LoTV) (see Bertsekas & Tsitsiklis, 2002), which gives

$$\begin{aligned} \text{Var} [\mathcal{B}(\pi_M(a^*|s, t))] &= \mathbb{E}_{h_s \sim p_t^{\pi_{\text{NM}}}} [\text{Var} [\mathcal{B}(\pi_{\text{NM}}(a^*|hs))] \\ &\quad + \text{Var}_{h_s \sim p_t^{\pi_{\text{NM}}}} [\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^*|hs))]], \quad \forall s \in \mathcal{S}. \end{aligned}$$

Then, exploiting the determinism of  $\pi_{\text{NM}}$  (through Lemma 4.4), it is straightforward to see that  $\mathbb{E}_{h_s \sim p_t^{\pi_{\text{NM}}}} [\text{Var} [\mathcal{B}(\pi_{\text{NM}}(a^*|hs))] = 0$ , which concludes the proof.<sup>3</sup>  $\square$

Unsurprisingly, Lemma 4.5 shows that, whenever the optimal strategy for (2) (i.e., the non-Markovian  $\pi_{\text{NM}}$ ) requires to adapt its decision in a state  $s$  according to the history that led to it ( $hs$ ), an optimal Markovian policy for the same objective (i.e.,  $\pi_M$ ) must necessarily be randomized. This is crucial to prove the following result, which establishes lower and upper bounds  $\underline{\mathcal{R}}_{T-t}, \overline{\mathcal{R}}_{T-t}$  to the expected regret-to-go of any Markovian policy that optimizes (2).

**Lemma 4.6.** *Let  $\pi_M$  be an optimal Markovian policy  $\pi_M \in \arg \max_{\pi \in \Pi_M} \mathcal{E}(\pi)$  on a CMP  $\mathcal{M}$ . For any  $h_t \in \mathcal{H}_{[T]}$ , it holds  $\underline{\mathcal{R}}_{T-t}(\pi_M) \leq \mathcal{R}_{T-t}(\pi_M) \leq \overline{\mathcal{R}}_{T-t}(\pi_M)$  such that*

$$\begin{aligned} \underline{\mathcal{R}}_{T-t}(\pi_M) &= \frac{H^* - H_2^*}{\pi_M(a^*|s_t)} \text{Var}_{h_s \sim p_t^{\pi_{\text{NM}}}} [\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^*|hs_t))]], \\ \overline{\mathcal{R}}_{T-t}(\pi_M) &= \frac{H^* - H_*}{\pi_M(a^*|s_t)} \text{Var}_{h_s \sim p_t^{\pi_{\text{NM}}}} [\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^*|hs_t))]], \end{aligned}$$

where  $\pi_{\text{NM}} \in \arg \max_{\pi \in \Pi_{\text{NM}}^{\text{D}}} \mathcal{E}(\pi)$ , and  $H_*, H_2^*$  are given

<sup>3</sup>Note that the determinism of  $\pi_{\text{NM}}$  does not also imply  $\text{Var}_{h_s \sim p_t^{\pi_{\text{NM}}}} [\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^*|hs))] = 0$ , as the optimal action  $\bar{a} = \pi_{\text{NM}}(hs)$  may vary for different histories, which results in the inner expectations  $\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^*|hs))]$  being either 1 or 0.

by

$$\begin{aligned} H_* &= \min_{h \in \mathcal{H}_{T-t}} H(d_{h_t \oplus h}(\cdot)), \\ H_2^* &= \max_{h \in \mathcal{H}_{T-t} \setminus \mathcal{H}_{T-t}^*} H(d_{h_t \oplus h}(\cdot)) \\ \text{s.t. } \mathcal{H}_{T-t}^* &= \arg \max_{h \in \mathcal{H}_{T-t}} H(d_{h_t \oplus h}(\cdot)). \end{aligned}$$

*Proof Sketch.* The crucial idea to derive lower and upper bounds to the regret-to-go is to consider the impact of a sub-optimal action in the best-case and the worst-case CMP respectively (see Lemma B.2, B.1). This gives  $\mathcal{R}_{T-t}(\pi_M) \geq H^* - \pi_M(a^*|s_t)H^* - (1 - \pi_M(a^*|s_t))H_2^*$  and  $\mathcal{R}_{T-t}(\pi_M) \leq H^* - \pi_M(a^*|s_t)H^* - (1 - \pi_M(a^*|s_t))H_*$ . Then, with Lemma 4.5 we get  $\text{Var}[\mathcal{B}(\pi_M(a^*|s_t))] = \pi_M(a^*|s_t)(1 - \pi_M(a^*|s_t)) = \text{Var}_{hs \sim p_t^{\pi_{NM}}}[\mathbb{E}[\mathcal{B}(\pi_{NM}(a^*|hs_t))]]$ , which concludes the proof.  $\square$

Finally, the result in Theorem 4.2 is a direct consequence of Lemma 4.6. Note that the upper and lower bounds on the regret-to-go are strictly positive whenever  $\pi_M(a^*|s_t) < 1$ , as it is stated in Corollary 4.3.

## 5. Complexity Analysis

Having established the importance of non-Markovianity in dealing with MSE exploration in a finite-sample regime, it is worth considering how hard it is to optimize the objective 2 within the class of non-Markovian policies. Especially, we aim at characterizing the complexity of the problem:

$$\Psi_0 := \underset{\pi \in \Pi_{NM}}{\text{maximize}} \mathcal{E}(\pi),$$

defined over a CMP  $\mathcal{M}$ . Before going into the details of the analysis, we provide a couple of useful definitions for the remainder of the section, whereas we leave to (Arora & Barak, 2009) an extended review of complexity theory.

**Definition 5.1** (Many-to-one Reductions). *We denote as  $A \leq_m B$  a many-to-one reduction from  $A$  to  $B$ .*

**Definition 5.2** (Polynomial Reductions). *We denote as  $A \leq_p B$  a polynomial-time (Turing) reduction from  $A$  to  $B$ .*

Then, we recall that  $\Psi_0$  can be rewritten as the problem of finding a reward-maximizing Markovian policy, i.e.,  $\tilde{\pi}_M \in \arg \max_{\pi \in \Pi_M} \mathcal{J}_{\tilde{\mathcal{M}}_T^R}(\pi)$ , over a convenient extended MDP  $\tilde{\mathcal{M}}_T^R$  obtained from CMP  $\mathcal{M}$  (see the proof of Lemma 4.4 for further details). We call this problem  $\tilde{\Psi}_0$  and we note that  $\tilde{\Psi}_0 \in \text{P}$ , as the problem of finding a reward-maximizing Markovian policy is well-known to be in P for any MDP (Papadimitriou & Tsitsiklis, 1987). However, the following lemma shows that it does not exist a many-to-one reduction from  $\Psi_0$  to  $\tilde{\Psi}_0$ .

**Lemma 5.3.** *A reduction  $\Psi_0 \leq_m \tilde{\Psi}_0$  does not exist.*

*Proof.* In general, coding any instance of  $\Psi_0$  in the representation required by  $\tilde{\Psi}_0$ , which is an extended MDP  $\tilde{\mathcal{M}}_T^R$ , holds exponential complexity w.r.t. the input of the initial instance of  $\Psi_0$ , i.e., a CMP  $\mathcal{M}$ . Indeed, to build the extended MDP  $\tilde{\mathcal{M}}_T^R$  from  $\mathcal{M}$ , we need to define the transition probabilities  $\tilde{P}(\tilde{s}'|\tilde{s}, \tilde{a})$  for every  $\tilde{s}' \in \tilde{\mathcal{S}}, \tilde{a} \in \tilde{\mathcal{A}}, \tilde{s} \in \tilde{\mathcal{S}}$ . Whereas the action space remains unchanged  $\tilde{\mathcal{A}} = \mathcal{A}$ , the extended state space  $\tilde{\mathcal{S}}$  has cardinality  $|\tilde{\mathcal{S}}| = S^T$  in general, which grows exponentially in  $T$ .  $\square$

The latter result informally suggests that  $\Psi_0 \notin \text{P}$ . Indeed, we can now prove the main theorem of this section, which shows that  $\Psi_0$  is NP-hard under the common assumption that  $\text{P} \neq \text{NP}$ .

**Theorem 5.4.**  *$\Psi_0$  is NP-hard.*

*Proof Sketch.* To prove the theorem, it is sufficient to show that there exists a problem  $\Psi_c \in \text{NP-hard}$  so that  $\Psi_c \leq_p \Psi_0$ . We show this by reducing 3SAT, which is a well-known NP-complete problem, to  $\Psi_0$ . To derive the reduction we consider two intermediate problems, namely  $\Psi_1$  and  $\Psi_2$ . Especially, we aim to show that the following chain of reductions holds

$$\Psi_0 \geq_m \Psi_1 \geq_p \Psi_2 \geq_p \text{3SAT}.$$

First, we define  $\Psi_1$  and we prove that  $\Psi_0 \geq_m \Psi_1$ . Informally,  $\Psi_1$  is the problem of finding a reward-maximizing Markovian policy  $\pi_M \in \Pi_M$  w.r.t. the entropy objective (2) encoded through a reward function in a convenient POMDP  $\tilde{\mathcal{M}}_\Omega^R$ . We can build  $\tilde{\mathcal{M}}_\Omega^R$  from the CMP  $\mathcal{M}$  similarly as the extended MDP  $\tilde{\mathcal{M}}_T^R$  (see Section 2 and the proof of Lemma 4.4 for details), except that the agent only access the observation space  $\tilde{\Omega}$  instead of the extended state space  $\tilde{\mathcal{S}}$ . In particular, we define  $\tilde{\Omega} = \mathcal{S}$  (note that  $\mathcal{S}$  is the state space of the original CMP  $\mathcal{M}$ ), and  $\tilde{O}(\tilde{o}|\tilde{s}) = \tilde{s}_{-1}$ . Then, the reduction  $\Psi_0 \geq_m \Psi_1$  works as follows. We denote as  $\mathcal{I}_{\Psi_i}$  the set of possible instances of problem  $\Psi_i$ . We show that  $\Psi_0$  is harder than  $\Psi_1$  by defining the polynomial-time functions  $\psi$  and  $\phi$  such that any instance of  $\Psi_1$  can be rewritten through  $\psi$  as an instance of  $\Psi_0$ , and a solution  $\pi_{NM}^* \in \Pi_{NM}$  for  $\Psi_0$  can be converted through  $\phi$  into a solution  $\pi_M^* \in \Pi_M$  for the original instance of  $\Psi_1$ . The function  $\psi$  sets  $\mathcal{S} = \tilde{\Omega}$  and derives the transition model of  $\mathcal{M}$  from the one of  $\tilde{\mathcal{M}}_\Omega^R$ , while  $\phi$  converts the optimal solution of  $\Psi_0$  by computing  $\pi_M^*(a|o, t) = \sum_{ho \in \mathcal{H}_o} p_T^{\pi_{NM}^*}(ho) \pi_{NM}^*(a|ho)$ , where  $\mathcal{H}_o$  stands for the set of histories  $h \in \mathcal{H}_t$  ending in the observation  $o \in \Omega$ . Thus, we have that  $\Psi_0 \geq_m \Psi_1$  holds. We now define  $\Psi_2$  as the policy existence problem w.r.t. the problem statement of  $\Psi_1$ . Hence,  $\Psi_2$  is the problem of determining whether the value of a reward-maximizing

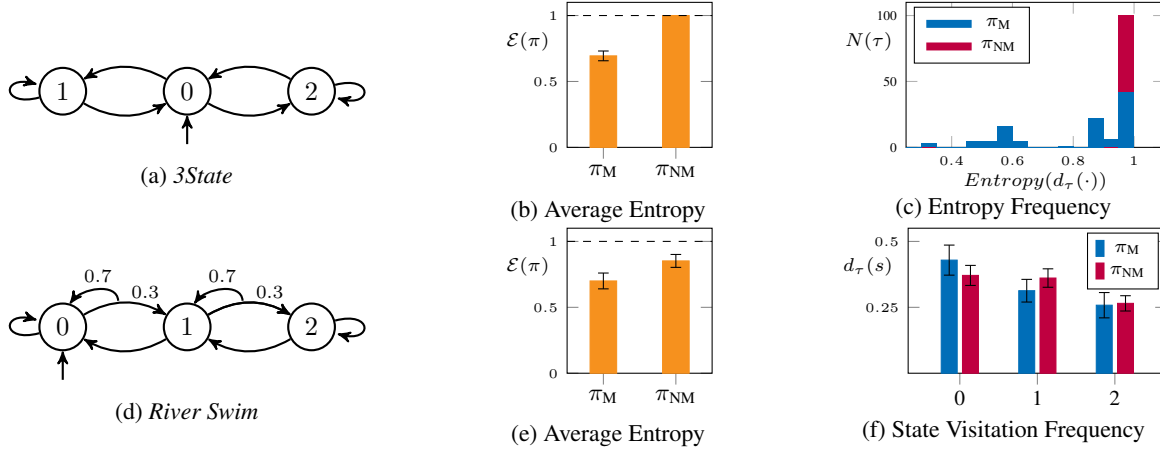


Figure 2. In (a, d), we illustrate the *3State* and *River Swim* CMPs. Then, we report the average entropy induced by an optimal (stationary) Markovian policy  $\pi_M$  and an optimal non-Markovian policy  $\pi_{NM}$  in the *3State* ( $T = 9$ ) (b) and the *River Swim* ( $T = 10$ ) (e). In (c) we report the entropy frequency in the *3State*, in (f) the state visitation frequency in the *River Swim*. We provide 95% c.i. over 100 runs.

Markovian policy  $\pi_M^* \in \arg \max_{\pi \in \Pi_M} \mathcal{J}_{\mathcal{M}_\Omega^R}(\pi)$  is greater than 0. Since computing an optimal policy in POMDPs is in general harder than the relative policy existence problem (Lusena et al., 2001, Section 3), we have that  $\Psi_1 \geq_p \Psi_2$ . For the last reduction, i.e.,  $\Psi_2 \geq_p 3SAT$ , we extend the proof of Theorem 4.13 in (Mundhenk et al., 2000), which states that the policy existence problem for POMDPs is NP-complete. In particular, we show that this holds within the restricted class of POMDPs defined in  $\Psi_1$ . Since the chain  $\Psi_0 \geq_m \Psi_1 \geq_p \Psi_2 \geq_p 3SAT$  holds, we have that  $\Psi_0 \geq_p 3SAT$ . Since  $3SAT \in NP$ -complete, we can conclude that  $\Psi_0$  is NP-hard.  $\square$

Having established the hardness of the optimization of  $\Psi_0$ , one could now question whether the problem  $\Psi_0$  is instead easy to verify ( $\Psi_0 \in NP$ ), from which we would conclude that  $\Psi_0 \in NP$ -complete. Whereas we doubt that this problem is significantly easier to verify than to optimize, the focus of this work is on its optimization version, and we thus leave as future work a finer analysis to show that  $\Psi_0 \notin NP$ .

## 6. Numerical Validation

Despite the hardness result of Theorem 5.4, we provide a brief numerical validation around the potential of non-Markovianity in MSE exploration. Crucially, the reported analysis is limited to simple domains and short time horizons, and it has to be intended as an illustration of the theoretical claims reported in previous sections. For the sake of simplicity, in this analysis we consider stationary policies for the Markovian set, though similar results can be obtained for time-variant strategies as well (in stochastic environments). Whereas a comprehensive evaluation of the practical benefits of non-Markovianity in MSE exploration is left as future work, we discuss in Section 7 why

we believe that the development of scalable methods is not hopeless even in this challenging setting.

In this section, we consider a *3State* ( $S = 3, A = 2, T = 9$ ), which is a simple abstraction of the two-rooms in Figure 1, and a *River Swim* (Strehl & Littman, 2008) ( $S = 3, A = 2, T = 10$ ) that are depicted in Figure 2a, 2d respectively. Especially, we compare the expected entropy (2) achieved by an optimal non-Markovian policy  $\pi_{NM} \in \arg \max_{\pi \in \Pi_{NM}} \mathcal{E}(\pi)$ , which is obtained by solving the extended MDP as described in the proof of Lemma 4.4, against an optimal Markovian policy  $\pi_M \in \arg \max_{\pi \in \Pi_M} \mathcal{E}(\pi)$ . In confirmation of the result in Theorem 4.2,  $\pi_M$  cannot match the performance of  $\pi_{NM}$  (see Figure 2b, 2e). In *3State*, an optimal strategy requires going left when arriving in state 0 from state 2 and vice versa. The policy  $\pi_{NM}$  is able to do that, and it always realizes the optimal trajectory (Figure 2c). Instead,  $\pi_M$  is uniform in 0 and it often runs into sub-optimal trajectories. In the *River Swim*, the main hurdle is to reach state 2 from the initial one. Whereas  $\pi_M$  and  $\pi_{NM}$  are equivalently good in doing so, as reported in Figure 2f, only the non-Markovian strategy is able to balance the visitations in the previous states when it eventually reaches 2. The difference is already noticeable with a short horizon and it would further increase with a longer  $T$ .

## 7. Discussion and Conclusion

In the previous sections, we detailed the importance of non-Markovianity when optimizing a finite-sample MSE objective, but we also proved that the corresponding optimization problem is NP-hard in its general formulation. Despite the hardness result, we believe that it is not hopeless to learn exploration policies with some form of non-Markovianity, while still preserving an edge over Markovian strategies.

In the following paragraphs, we discuss potential avenues to derive practical methods for relevant relaxations to the general class of non-Markovian policies.

**Finite-Length Histories** Throughout the paper, we considered non-Markovian policies that condition their decisions on histories of arbitrary length, i.e.,  $\pi : \mathcal{H} \rightarrow \Delta(\mathcal{A})$ . However, the complexity of optimizing such policies grows exponentially with the length of the history. To avoid this exponential blowup, one can define a class of non-Markovian policies  $\pi : \mathcal{H}_H \rightarrow \Delta(\mathcal{A})$  in which the decisions are conditioned on histories of a finite length  $H > 1$  that are obtained from a sliding window on the full history. The optimal policy within this class would still retain better regret guarantees than an optimal Markovian policy, but it would not achieve zero regret in general. With the length parameter  $H$  one can trade-off the learning complexity with the regret according to the structure of the domain. For instance,  $H = 2$  would be sufficient to achieve zero regret in the *3State* domain, whereas in the *River Swim* domain any  $H < T$  would cause some positive regret.

**Compact Representations of the History** Instead of setting a finite length  $H$ , one can choose to perform function approximation on the full history to obtain a class of policies  $\pi : f(\mathcal{H}) \rightarrow \Delta(\mathcal{A})$ , where  $f$  is a function that maps an history  $h$  to some compact representation. An interesting option is to use the notion of *eligibility traces* (Sutton & Barto, 2018) to encode the information of  $h$  in a vector of length  $S$ , which is updated as  $z_{t+1} \leftarrow \lambda z_t + \mathbf{1}_{s_t}$ , where  $\lambda \in (0, 1)$  is a discount factor,  $\mathbf{1}_{s_t}$  is a vector with a unit entry at the index  $s_t$ , and  $z_0 = 0$ . The discount factor  $\lambda$  acts as a smoothed version of the length parameter  $H$ , and it can be dynamically adapted while learning. Indeed, this eligibility traces representation is particularly convenient for policy optimization (Deisenroth et al., 2013), in which we could optimize in turn a parametric policy over actions  $\pi_{\theta}(\cdot | z, \lambda)$  and a parametric policy over the discount  $\pi_{\nu}(\lambda)$ . To avoid a direct dependence on  $S$ , one can define the vector  $z$  over a discretization of the state space.

**Deep Recurrent Policies** Another noteworthy way to do function approximation on the history is to employ recurrent neural networks (Williams & Zipser, 1989; Hochreiter & Schmidhuber, 1997) to represent the non-Markovian policy. This kind of recurrent architecture is already popular in RL. In this paper we are providing the theoretical ground to motivate the use of deep recurrent policies to address maximum state entropy exploration.

**Non-Markovian Control with Tree Search** In principle, one can get a realization of actions from the optimal non-Markovian policy without ever computing it, e.g., by employing a Monte-Carlo Tree Search (MCTS) (Kocsis & Szepesvári, 2006) approach to select the next action to take. Given the current state  $s_t$  as a root, we can build the tree

of trajectories from the root through repeated simulations of potential action sequences. With a sufficient number of simulations and a sufficiently deep tree, we are guaranteed to select the optimal action at the root. If the horizon is too long, we can still cut the tree at any depth and approximately evaluate a leaf node with the entropy induced by the path from the root to the leaf. The drawback of this procedure is that we require to access a simulator with reset (or a reliable estimate of the transition model) to actually build the tree.

Having reported interesting directions to learn non-Markovian exploration policies in practice, we would like to mention some relevant online RL settings that might benefit from such exploration policies. We leave as future work a formal definition of the settings and an empirical study.

**Single-Trial RL** In many relevant real-world scenarios, where data collection might be costly or non-episodic in nature, we cannot afford multiple trials to achieve the desired exploration of the environment. Non-Markovian exploration policies guarantee a good coverage of the environment in a single trial and they are particularly suitable for online learning processes.

**Learning in Latent MDPs** In a latent MDP scenario (Halak et al., 2015; Kwon et al., 2021) an agent interacts with an (unknown) environment drawn from a class of MDPs to solve an online RL task. A non-Markovian exploration policy pre-trained on the whole class could exploit the memory to perform a fast identification of the specific context that has been drawn, quickly adapting to the optimal environment-specific policy.

In this paper we focus on the gap between non-Markovian and Markovian policies, which can be either stationary or time-variant. Future works might consider the role of stationarity (see also Akshay et al., 2013; Laroche et al., 2022), such as establishing under which conditions stationary strategies are sufficient in this setting. Finally, here we focus on state distributions, which is most common in the MSE literature, but similar results could be extended to state-action distributions with minor modifications.

To conclude, we believe that this work sheds some light on the, previously neglected, importance of non-Markovianity to address maximum state entropy exploration. Although it brings a negative result about the computational complexity of the problem, we believe it can provide inspiration for future empirical and theoretical contributions on the matter.

## Acknowledgements

We would like to thank the reviewers of this paper for their feedbacks and useful advices. We also thank Romain Laroche and Remi Tachet des Combes for having signalled a technical error in a previous draft of the manuscript.



## References

- Akshay, S., Bertrand, N., Haddad, S., and Helouet, L. The steady-state control problem for Markov decision processes. In *International Conference on Quantitative Evaluation of Systems*, 2013.
- Arora, S. and Barak, B. *Computational complexity: a modern approach*. Cambridge University Press, 2009.
- Astrom, K. J. Optimal control of Markov decision processes with incomplete state estimation. *Journal Mathematical Analysis and Applications*, 1965.
- Bertsekas, D. P. and Tsitsiklis, J. N. *Introduction to probability*. Athena Scientific Belmont, MA, 2002.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 2020.
- Campos, V., Sprechmann, P., Hansen, S., Barreto, A., Kapurowski, S., Vitvitskyi, A., Badia, A. P., and Blundell, C. Coverage as a principle for discovering transferable behavior in reinforcement learning. *arXiv preprint arXiv:2102.13515*, 2021.
- Cheung, W. C. Exploration-exploitation trade-off in reinforcement learning on online Markov decision processes with global concave rewards. *arXiv preprint arXiv:1905.06466*, 2019a.
- Cheung, W. C. Regret minimization for reinforcement learning with vectorial feedback and complex objectives. In *Advances in Neural Information Processing Systems*, 2019b.
- Deisenroth, M. P., Neumann, G., Peters, J., et al. A survey on policy search for robotics. *Foundations and Trends in Robotics*, 2013.
- Erhan, D., Manzagol, P.-A., Bengio, Y., Bengio, S., and Vincent, P. The difficulty of training deep architectures and the effect of unsupervised pre-training. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2009.
- Erhan, D., Courville, A., Bengio, Y., and Vincent, P. Why does unsupervised pre-training help deep learning? In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2010.
- Guo, Z. D., Azar, M. G., Saade, A., Thakoor, S., Piot, B., Pires, B. A., Valko, M., Mesnard, T., Lattimore, T., and Munos, R. Geometric entropic exploration. *arXiv preprint arXiv:2101.02055*, 2021.
- Hallak, A., Di Castro, D., and Mannor, S. Contextual Markov decision processes. *arXiv preprint arXiv:1502.02259*, 2015.
- Hazan, E., Kakade, S., Singh, K., and Van Soest, A. Provably efficient maximum entropy exploration. In *Proceedings of the International Conference on Machine Learning*, 2019.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Computation*, 1997.
- Jin, C., Krishnamurthy, A., Simchowitz, M., and Yu, T. Reward-free exploration for reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, 2020.
- Kaelbling, L. P., Littman, M. L., and Cassandra, A. R. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 1998.
- Kocsis, L. and Szepesvári, C. Bandit based Monte-Carlo planning. In *European Conference on Machine Learning*, 2006.
- Kwon, J., Efroni, Y., Caramanis, C., and Mannor, S. RL for latent MDPs: Regret guarantees and a lower bound. In *Advances in Neural Information Processing Systems*, 2021.
- Laroche, R., Combes, R. T. d., and Buckman, J. Non-Markovian policies occupancy measures. *arXiv preprint arXiv:2205.13950*, 2022.
- Laskin, M., Yarats, D., Liu, H., Lee, K., Zhan, A., Lu, K., Cang, C., Pinto, L., and Abbeel, P. URLB: Unsupervised reinforcement learning benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.
- Lee, L., Eysenbach, B., Parisotto, E., Xing, E., Levine, S., and Salakhutdinov, R. Efficient exploration via state marginal matching. *arXiv preprint arXiv:1906.05274*, 2019.
- Liu, H. and Abbeel, P. APS: Active pretraining with successor features. In *Proceedings of the International Conference on Machine Learning*, 2021a.
- Liu, H. and Abbeel, P. Behavior from the void: Unsupervised active pre-training. In *Advances in Neural Information Processing Systems*, 2021b.

- Lusena, C., Goldsmith, J., and Mundhenk, M. Nonapproximability results for partially observable Markov decision processes. *Journal of Artificial Intelligence Research*, 2001.
- Mundhenk, M., Goldsmith, J., Lusena, C., and Allender, E. Complexity of finite-horizon Markov decision process problems. *Journal of the ACM (JACM)*, 2000.
- Mutti, M. and Restelli, M. An intrinsically-motivated approach for learning highly exploring and fast mixing policies. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- Mutti, M., Pratisoli, L., and Restelli, M. Task-agnostic exploration via policy gradient of a non-parametric state entropy estimate. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- Mutti, M., Mancassola, M., and Restelli, M. Unsupervised reinforcement learning in multiple environments. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
- Nedergaard, A. and Cook, M. k-means maximum entropy exploration. *arXiv preprint arXiv:2205.15623*, 2022.
- Papadimitriou, C. H. and Tsitsiklis, J. N. The complexity of Markov decision processes. *Mathematics of Operations Research*, 1987.
- Peters, J. and Schaal, S. Reinforcement learning of motor skills with policy gradients. *Neural networks*, 2008.
- Puterman, M. L. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Seo, Y., Chen, L., Shin, J., Lee, H., Abbeel, P., and Lee, K. State entropy maximization with random encoders for efficient exploration. In *Proceedings of the International Conference on Machine Learning*, 2021.
- Strehl, A. L. and Littman, M. L. An analysis of model-based interval estimation for Markov decision processes. *Journal of Computer and System Sciences*, 2008.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Tarbouriech, J. and Lazaric, A. Active exploration in Markov decision processes. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2019.
- Tarbouriech, J., Pirota, M., Valko, M., and Lazaric, A. A provably efficient sample collection strategy for reinforcement learning. In *Advances in Neural Information Processing Systems*, 2021.
- Watkins, C. J. and Dayan, P. Q-learning. *Machine learning*, 1992.
- Williams, R. J. and Zipser, D. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1989.
- Yarats, D., Fergus, R., Lazaric, A., and Pinto, L. Reinforcement learning with prototypical representations. In *Proceedings of the International Conference on Machine Learning*, 2021.
- Zhang, C., Cai, Y., Huang, L., and Li, J. Exploration by maximizing Rényi entropy for reward-free RL framework. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- Zhang, J., Koppel, A., Bedi, A. S., Szepesvari, C., and Wang, M. Variational policy gradient method for reinforcement learning with general utilities. In *Advances in Neural Information Processing Systems*, 2020.

Table 1. Overview of the methods addressing MSE exploration in a controlled Markov process. For each method, we report the nature of the corresponding MSE objective, i.e., the entropy function (Entropy), whether it considers stationary, discounted, or marginal distributions (Distribution), and if it accounts for the state space  $\mathcal{S}$  or the state-action space  $\mathcal{S}A$  (Space). We also specify if the method learns a single policy rather than a mixture of policies (Mixture), and if it supports non-parametric entropy estimation (Non-Parametric).

Algorithm	Entropy	Distribution	Space	Mixture	Non-Parametric
MaxEnt (Hazan et al., 2019)	Shannon	discounted	state	✓	✗
FW-AME (Tarbouriech & Lazaric, 2019)	Shannon	stationary	state-action	✓	✗
SMM (Lee et al., 2019)	Shannon	marginal	state	✓	✗
IDE <sup>3</sup> AL (Mutti & Restelli, 2020)	Shannon	stationary	state-action	✗	✗
MEPOL (Mutti et al., 2021)	Shannon	marginal	state	✗	✓
MaxRényi (Zhang et al., 2021)	Rényi	discounted	state-action	✗	✗
GEM (Guo et al., 2021)	geometry-aware	marginal	state	✗	✗
APT (Liu & Abbeel, 2021b)	Shannon	marginal	state	✗	✓
RE3 (Seo et al., 2021)	Shannon	marginal	state	✗	✓
Proto-RL (Yarats et al., 2021)	Shannon	marginal	state	✗	✓
APS (Liu & Abbeel, 2021a)	Shannon	marginal	state	✓	✓

## A. Related Work

Hazan et al. (2019) were the first to consider an entropic measure over the state distribution as a sensible learning objective for an agent interacting with a reward-free environment (Jin et al., 2020). Especially, they propose an algorithm, called MaxEnt, that learns a mixture of policies that collectively maximize the Shannon entropy of the discounted state distribution, i.e., (1). The final mixture is learned through a conditional gradient method, in which the algorithm iteratively estimates the state distribution of the current mixture to define an intrinsic reward function, and then identifies the next policy to be added by solving a specific RL sub-problem with this reward. A similar methodology has been obtained by Lee et al. (2019) from a game-theoretic perspective on the MSE exploration problem. Their algorithm, called SMM, targets the Shannon entropy of the marginal state distribution instead of the discounted distribution of MaxEnt. Another approach based on the conditional gradient method is FW-AME (Tarbouriech & Lazaric, 2019), which learns a mixture of policies to maximize the entropy of the stationary state-action distribution. As noted in (Tarbouriech & Lazaric, 2019), the mixture of policies might suffer a slow mixing to the asymptotic distribution for which the entropy is maximized. In (Mutti & Restelli, 2020), the authors present a method (IDE<sup>3</sup>AL) to learn a single exploration policy that simultaneously accounts for the entropy of the stationary state-action distribution and the mixing time.

Even if they are sometimes evaluated on continuous domains (especially (Hazan et al., 2019; Lee et al., 2019)), the methods we mentioned require an accurate estimate of either the state distribution (Hazan et al., 2019; Lee et al., 2019) or the transition model (Tarbouriech & Lazaric, 2019; Mutti & Restelli, 2020), which hardly scales to high-dimensional domains. A subsequent work by Mutti et al. (2021) proposes an approach to estimate the entropy of the state distribution through a non-parametric method, and then to directly optimize the estimated entropy via policy optimization. Their algorithm, called MEPOL, is able to learn a single exploration policy that maximizes the entropy of the marginal state distribution in challenging continuous control domains. Liu & Abbeel (2021b) combine non-parametric entropy estimation with learned state representations into an algorithm, called APT, that successfully addresses MSE exploration problems in visual-inputs domains. Seo et al. (2021) shows that even random state representations are sufficient to learn MSE exploration policies from visual inputs. On a similar line, Yarats et al. (2021) consider simultaneously learning state representations and a basis for the latent space (or prototypical representations) to help reducing the variance of the entropy estimates. Finally, Liu & Abbeel (2021a) consider a method, called APS, to learn a set of code-conditioned policies that collectively maximizes the MSE objective by coupling non-parametric entropy estimation and successor representation.

Whereas all of the previous approaches accounts for the Shannon entropy in their objectives, recent works (Zhang et al., 2021; Guo et al., 2021) consider alternative formulations. Zhang et al. (2021) argues that the Rényi entropy provides a superior incentive to cover all of the corresponding space than the Shannon entropy, and they propose a method to optimize the Rényi of the state-action distribution via gradient ascent (MaxRényi). On an orthogonal direction, the authors of (Guo et al., 2021) consider a reformulation of the entropy function that accounts for the underlying geometry of the space. They present a method, called GEM, to learn an optimal policy for the geometry-aware entropy objective.

### A.1. Online Learning of Global Concave Rewards

Another interesting pair of related works (Cheung, 2019a;b) addresses a reinforcement learning problem for the maximization of concave functions of vectorial rewards, which in a special case (Cheung, 2019a, Section 6.2) is akin to our objective function (2). Beyond this similarity in the objective definition, those works and our paper differ for some crucial aspects, and the contributions are essentially non-overlapping. On the one hand, (Cheung, 2019a;b) deal with an online learning problem, in which they care for the performance of the policies deployed during the learning process, whereas we only consider the performance of the optimal policy. On the other hand, we aim to compare the classes of non-Markovian and Markovian policies respectively, whereas they consider non-stationary or adaptive strategies to maximize the online objective. Finally, their definition of regret is based on an infinite-samples relaxation of the problem, whereas we account for the performance of the optimal general policy w.r.t. the finite-sample objective (2) to define our regret-to-go (Definition 4.1).

## B. Missing Proofs

### B.1. Proofs of Section 3

**Theorem 3.1.** *Let  $x \in \{\infty, \gamma, T\}$ , and let  $\mathcal{D}_{\text{NM}}^x = \{d_x^\pi(\cdot) : \pi \in \Pi_{\text{NM}}\}$ ,  $\mathcal{D}_{\text{M}}^x = \{d_x^\pi(\cdot) : \pi \in \Pi_{\text{M}}\}$  the corresponding sets of state distributions over a CMP. We can prove that:*

- (i) *The sets of stationary state distributions are equivalent  $\mathcal{D}_{\text{NM}}^\infty \equiv \mathcal{D}_{\text{M}}^\infty$ ;*
- (ii) *The sets of discounted state distributions are equivalent  $\mathcal{D}_{\text{NM}}^\gamma \equiv \mathcal{D}_{\text{M}}^\gamma$  for any  $\gamma$ ;*
- (iii) *The sets of marginal state distributions are equivalent  $\mathcal{D}_{\text{NM}}^T \equiv \mathcal{D}_{\text{M}}^T$  for any  $T$ .*

*Proof.* First, note that a non-Markovian policy  $\pi \in \Pi_{\text{NM}}$  can always reduce to a Markovian policy  $\pi \in \Pi_{\text{M}}$  by conditioning the decision rules on the history length. Thus,  $\mathcal{D}_{\text{NM}}^x \supseteq \mathcal{D}_{\text{M}}^x$  is straightforward for any  $x \in \{\infty, \gamma, T\}$ . From the derivations in (Puterman, 2014, Theorem 5.5.1), we have that  $\mathcal{D}_{\text{M}}^x \supseteq \mathcal{D}_{\text{NM}}^x$  as well. Indeed, for any non-Markovian policy  $\pi \in \Pi_{\text{NM}}$ , we can build a (non-stationary) Markovian policy  $\pi' \in \Pi_{\text{M}}$  as

$$\pi' = (\pi'_1, \pi'_2, \dots, \pi'_t, \dots), \quad \text{such that } \pi'_t(a|s) = \frac{d_t^\pi(s, a)}{d_t^\pi(s)}, \quad \forall s \in \mathcal{S}, \forall a \in \mathcal{A}.$$

For  $t = 0$ , we have that  $d_0^\pi(\cdot) = d_0^{\pi'}(\cdot) = \mu(\cdot)$ , which is the initial state distribution. We proceed by induction to show that if  $d_{t-1}^\pi(\cdot) = d_{t-1}^{\pi'}(\cdot)$ , then we have

$$\begin{aligned} d_t^{\pi'}(s) &= \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{t-1}^{\pi'}(s') \pi'_{t-1}(a|s') P(s|s', a) \\ &= \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} \frac{d_{t-1}^{\pi'}(s')}{d_{t-1}^{\pi'}(s')} d_{t-1}^{\pi'}(s', a) P(s|s', a) \\ &= \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{t-1}^\pi(s', a) P(s|s', a) \\ &= d_t^\pi(s). \end{aligned}$$

Since  $d_t^\pi(s) = d_t^{\pi'}(s)$  holds for any  $t \geq 0$  and  $\forall s \in \mathcal{S}$ , we have  $d_\infty^\pi(\cdot) = d_\infty^{\pi'}(\cdot)$ ,  $d_\gamma^\pi(\cdot) = d_\gamma^{\pi'}(\cdot)$ ,  $d_T^\pi(\cdot) = d_T^{\pi'}(\cdot)$ , and thus  $\mathcal{D}_{\text{M}}^x \supseteq \mathcal{D}_{\text{NM}}^x$ . Then,  $\mathcal{D}_{\text{NM}}^x \equiv \mathcal{D}_{\text{M}}^x$  follows.  $\square$

**Corollary 3.2.** *For every CMP, there exists a Markovian policy  $\pi^* \in \Pi_{\text{M}}$  such that  $\pi^* \in \arg \max_{\pi \in \Pi} \mathcal{E}_\infty(\pi)$ .*

*Proof.* The result is straightforward from Theorem 3.1 and noting that the set of non-Markovian policies  $\Pi_{\text{NM}}$  with arbitrary history-length is as powerful as the general set of policies  $\Pi$ . Thus, for every policy  $\pi \in \Pi$  there exists a (possibly randomized) policy  $\pi' \in \Pi_{\text{M}}$  inducing the same (stationary, discounted or marginal) state distribution of  $\pi$ , i.e.,  $d^\pi(\cdot) = d^{\pi'}(\cdot)$ , which implies  $H(d^\pi(\cdot)) = H(d^{\pi'}(\cdot))$ . If it holds for any  $\pi \in \Pi$ , then it holds for  $\pi^* \in \arg \max_{\pi \in \Pi} H(d^\pi(\cdot))$ .  $\square$



## B.2. Proofs of Section 4

**Theorem 4.2** (Non-Markovian Optimality). *For every CMP  $\mathcal{M}$  and trajectory  $h_t \in \mathcal{H}_{[T]}$ , there exists a deterministic non-Markovian policy  $\pi_{\text{NM}} \in \Pi_{\text{NM}}^{\text{D}}$  that suffers zero regret-to-go  $\mathcal{R}_{T-t}(\pi_{\text{NM}}, h_t) = 0$ , whereas for any  $\pi_{\text{M}} \in \Pi_{\text{M}}$  we have  $\mathcal{R}_{T-t}(\pi_{\text{M}}, h_t) \geq 0$ .*

*Proof.* The result  $\mathcal{R}_{T-t}(\pi_{\text{NM}}, h_t) = 0$  for a policy  $\pi_{\text{NM}} \in \Pi_{\text{D}}^{\text{NM}}$  is a direct implication of Lemma 4.4, whereas  $\mathcal{R}_{T-t}(\pi_{\text{M}}, h_t) \geq 0$  for any  $\pi_{\text{M}} \in \Pi_{\text{M}}$  is given by Lemma 4.6, which states that even an optimal Markovian policy  $\pi_{\text{M}}^* \in \arg \max_{\pi \in \Pi_{\text{M}}} \mathcal{E}(\pi)$  suffers expected regret-to-go  $\mathcal{R}_{T-t}(\pi_{\text{M}}^*) \geq 0$ .  $\square$

**Corollary 4.3** (Sufficient Condition). *For every CMP  $\mathcal{M}$  and trajectory  $h_t \in \mathcal{H}_{[T]}$  for which any optimal Markovian policy  $\pi_{\text{M}} \in \Pi_{\text{M}}$  is randomized (i.e., stochastic) in  $s_t$ , we have strictly positive regret-to-go  $\mathcal{R}_{T-t}(\pi_{\text{M}}, h_t) > 0$ .*

*Proof.* This result is a direct consequence of the combination of Lemma 4.5 and Lemma 4.6. Indeed, if the policy  $\pi_{\text{M}} \in \arg \max_{\pi \in \Pi_{\text{M}}} \mathcal{E}(\pi)$  is randomized in  $s_t$  we have

$$0 < \text{Var} [\mathcal{B}(\pi_{\text{M}}(a^* | s_t))] = \text{Var}_{h s_t \sim p_t^{\pi_{\text{NM}}}} [\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^* | h s_t))]],$$

from Lemma 4.5, which gives a lower bound to the expected regret-to-go  $\mathcal{R}_{T-t}(\pi_{\text{M}}, h_t) > 0$  through Lemma 4.6.  $\square$

**Lemma 4.5.** *Let  $\pi_{\text{NM}} \in \Pi_{\text{NM}}^{\text{D}}$  a non-Markovian policy such that  $\pi_{\text{NM}} \in \arg \max_{\pi \in \Pi} \mathcal{E}(\pi)$  on a CMP  $\mathcal{M}$ . For a fixed history  $h_t \in \mathcal{H}_t$  ending in state  $s$ , the variance of the event of an optimal Markovian policy  $\pi_{\text{M}} \in \arg \max_{\pi \in \Pi_{\text{M}}} \mathcal{E}(\pi)$  taking  $a^* = \pi_{\text{NM}}(h_t)$  in  $s$  is given by*

$$\text{Var} [\mathcal{B}(\pi_{\text{M}}(a^* | s, t))] = \text{Var}_{h s \sim p_t^{\pi_{\text{NM}}}} [\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^* | h s))]],$$

where  $h s \in \mathcal{H}_t$  is any history of length  $t$  such that the final state is  $s$ , i.e.,  $h s := (h_{t-1} \in \mathcal{H}_{t-1}) \oplus s$ , and  $\mathcal{B}(x)$  is a Bernoulli with parameter  $x$ .

*Proof.* Let us consider the random variable  $A \sim \mathcal{P}$  denoting the event “the agent takes action  $a^* \in \mathcal{A}$ ”. Through the law of total variance (Bertsekas & Tsitsiklis, 2002), we can write the variance of  $A$  given  $s \in \mathcal{S}$  and  $t \geq 0$  as

$$\begin{aligned} \text{Var} [A | s, t] &= \mathbb{E} [A^2 | s, t] - \mathbb{E} [A | s, t]^2 \\ &= \mathbb{E}_h [\mathbb{E} [A^2 | s, t, h]] - \mathbb{E}_h [\mathbb{E} [A | s, t, h]]^2 \\ &= \mathbb{E}_h [\text{Var} [A | s, t, h] + \mathbb{E} [A | s, t, h]^2] - \mathbb{E}_h [\mathbb{E} [A | s, t, h]]^2 \\ &= \mathbb{E}_h [\text{Var} [A | s, t, h]] + \mathbb{E}_h [\mathbb{E} [A | s, t, h]^2] - \mathbb{E}_h [\mathbb{E} [A | s, t, h]]^2 \\ &= \mathbb{E}_h [\text{Var} [A | s, t, h]] + \text{Var}_h [\mathbb{E} [A | s, t, h]]. \end{aligned} \quad (3)$$

Now let the conditioning event  $h$  be distributed as  $h \sim p_{t-1}^{\pi_{\text{NM}}}$ , so that the condition  $s, t, h$  becomes  $h s$  where  $h s = (s_0, a_0, s_1, \dots, s_t = s) \in \mathcal{H}_t$ , and let the variable  $A$  be distributed according to  $\mathcal{P}$  that maximizes the objective (2) given the conditioning. Hence, we have that the variable  $A$  on the left hand side of (3) is distributed as a Bernoulli  $\mathcal{B}(\pi_{\text{M}}(a^* | s, t))$ , where  $\pi_{\text{M}} \in \arg \max_{\pi \in \Pi_{\text{M}}} \mathcal{E}(\pi)$ , and the variable  $A$  on the right hand side of (4) is distributed as a Bernoulli  $\mathcal{B}(\pi_{\text{NM}}(a^* | h s))$ , where  $\pi_{\text{NM}} \in \arg \max_{\pi \in \Pi_{\text{NM}}} \mathcal{E}(\pi)$ . Thus, we obtain

$$\text{Var} [\mathcal{B}(\pi_{\text{M}}(a^* | s, t))] = \mathbb{E}_{h s \sim p_t^{\pi_{\text{NM}}}} [\text{Var} [\mathcal{B}(\pi_{\text{NM}}(a^* | h s))]] + \text{Var}_{h s \sim p_t^{\pi_{\text{NM}}}} [\mathbb{E} [\mathcal{B}(\pi_{\text{NM}}(a^* | h s))]]. \quad (4)$$

Under Assumption 1, we know from Lemma 4.4 that the policy  $\pi_{\text{NM}}$  is deterministic, i.e.,  $\pi_{\text{NM}} \in \Pi_{\text{D}}^{\text{NM}}$ , so that  $\text{Var} [\mathcal{B}(\pi_{\text{NM}}(a^* | h s))] = 0$  for every  $h s$ , which concludes the proof.  $\square$

**Lemma 4.6.** Let  $\pi_M$  be an optimal Markovian policy  $\pi_M \in \arg \max_{\pi \in \Pi_M} \mathcal{E}(\pi)$  on a CMP  $\mathcal{M}$ . For any  $h_t \in \mathcal{H}_{[T]}$ , it holds  $\underline{\mathcal{R}}_{T-t}(\pi_M) \leq \mathcal{R}_{T-t}(\pi_M) \leq \overline{\mathcal{R}}_{T-t}(\pi_M)$  such that

$$\begin{aligned}\underline{\mathcal{R}}_{T-t}(\pi_M) &= \frac{H^* - H_2^*}{\pi_M(a^*|s_t)} \text{Var}_{h_{s_t} \sim p_t^{\pi_{NM}}} \left[ \mathbb{E} [\mathcal{B}(\pi_{NM}(a^*|hs_t))] \right], \\ \overline{\mathcal{R}}_{T-t}(\pi_M) &= \frac{H^* - H_*}{\pi_M(a^*|s_t)} \text{Var}_{h_{s_t} \sim p_t^{\pi_{NM}}} \left[ \mathbb{E} [\mathcal{B}(\pi_{NM}(a^*|hs_t))] \right],\end{aligned}$$

where  $\pi_{NM} \in \arg \max_{\pi \in \Pi_{NM}^D} \mathcal{E}(\pi)$ , and  $H_*$ ,  $H_2^*$  are given by

$$\begin{aligned}H_* &= \min_{h \in \mathcal{H}_{T-t}} H(d_{h_t \oplus h}(\cdot)), \\ H_2^* &= \max_{h \in \mathcal{H}_{T-t} \setminus \mathcal{H}_{T-t}^*} H(d_{h_t \oplus h}(\cdot)) \\ \text{s.t. } \mathcal{H}_{T-t}^* &= \arg \max_{h \in \mathcal{H}_{T-t}} H(d_{h_t \oplus h}(\cdot)).\end{aligned}$$

*Proof.* From the definition of the expected regret-to-go (Definition 4.1), we have that

$$\mathcal{R}_{T-t}(\pi_M, h_t) = H^* - \mathbb{E}_{h_{T-t} \sim p_{T-t}^{\pi_M}} [H(d_{h_t \oplus h_{T-t}}(\cdot))],$$

in which we will omit  $h_t$  in the regret-to-go  $\mathcal{R}_{T-t}(\pi_M, h_t) = \mathcal{R}_{T-t}(\pi_M)$  as  $h_t$  is fixed by the statement. To derive a lower bound and an upper bound to  $\mathcal{R}_{T-t}(\pi_M)$  we consider the impact that taking a sub-optimal action  $a \in \mathcal{A} \setminus \{a^*\}$  in state  $s_t$  would have in a best-case and a worst-case CMP respectively, which is detailed in Lemma B.2 and Lemma B.1. Especially, we can write

$$\begin{aligned}\mathcal{R}_{T-t}(\pi_M) &= H^* - \mathbb{E}_{h_{T-t} \sim p_{T-t}^{\pi_M}} [H(d_{h_t \oplus h_{T-t}}(\cdot))] \\ &\geq H^* - \pi_M(a^*|s_t)H^* - (1 - \pi_M(a^*|s_t))H_2^* \\ &= (H^* - H_2^*)(1 - \pi_M(a^*|s_t))\end{aligned}$$

and

$$\begin{aligned}\mathcal{R}_{T-t}(\pi_M) &= H^* - \mathbb{E}_{h_{T-t} \sim p_{T-t}^{\pi_M}} [H(d_{h_t \oplus h_{T-t}}(\cdot))] \\ &\leq H^* - \pi_M(a^*|s_t)H^* - (1 - \pi_M(a^*|s_t))H_* \\ &= (H^* - H_*)(1 - \pi_M(a^*|s_t)).\end{aligned}$$

Then, we note that the event of taking a sub-optimal action  $a \in \mathcal{A} \setminus \{a^*\}$  with a policy  $\pi_M$  can be modelled by a Bernoulli distribution  $\mathcal{B}$  with parameter  $(1 - \pi_M(a^*|s_t))$ . By combining the equation of the variance of a Bernoulli random variable with Lemma 4.5 we obtain

$$\text{Var} [\mathcal{B}(\pi_M(a^*|s_t))] = \pi_M(a^*|s_t)(1 - \pi_M(a^*|s_t)) = \text{Var}_{hs \sim p_t^{\pi_{NM}}} \left[ \mathbb{E} [\mathcal{B}(\pi_{NM}(a^*|hs_t))] \right]$$

which gives

$$\begin{aligned}\mathcal{R}_{T-t}(\pi_M) &\geq \frac{H^* - H_2^*}{\pi_M(a^*|s_t)} \text{Var}_{hs \sim p_t^{\pi_{NM}}} \left[ \mathbb{E} [\mathcal{B}(\pi_{NM}(a^*|hs_t))] \right] := \underline{\mathcal{R}}_{T-t}(\pi_M) \\ \mathcal{R}_{T-t}(\pi_M) &\leq \frac{H^* - H_*}{\pi_M(a^*|s_t)} \text{Var}_{hs \sim p_t^{\pi_{NM}}} \left[ \mathbb{E} [\mathcal{B}(\pi_{NM}(a^*|hs_t))] \right] := \overline{\mathcal{R}}_{T-t}(\pi_M)\end{aligned}$$

□

**Lemma B.1** (Worst-Case CMP). *For any  $t$ -step trajectory  $h_t \in \mathcal{H}_{[T]}$ , taking a sub-optimal action  $a \in \mathcal{A} \setminus \{a^*\}$  at step  $t$  in the worst-case CMP  $\underline{\mathcal{M}}$  gives a final entropy*

$$H_* = \min_{h \in \mathcal{H}_{T-t}} H(d_{h_t \oplus h}(\cdot)),$$

where the minimum is attained by

$$\underline{h}_{T-t} = \left( s_i \in \arg \max_{s \in \mathcal{S}} d_{h_t}(s) \right)_{i=t+1}^T \in \arg \min_{h \in \mathcal{H}_{T-t}} H(d_{h_t \oplus h}(\cdot)).$$

*Proof.* The worst-case CMP  $\underline{\mathcal{M}}$  is designed such that the agent cannot recover from taking a sub-optimal action  $a_t \in \mathcal{A} \setminus \{a^*\}$  as it is absorbed by a worst-case state given the trajectory  $h_t$ . A worst-case state is one that maximizes the visitation frequency in  $h_t$ , i.e.,  $\underline{s} \in \arg \max_{s \in \mathcal{S}} d_{h_t}(s)$ , so that the visitation frequency becomes increasingly unbalanced. A sub-optimal action at the first step in  $\underline{\mathcal{M}}$  leads to  $T - 1$  visits to the initial state  $s_0 \sim \mu$ , and the final entropy is zero.  $\square$

**Lemma B.2** (Best-Case CMP). *For any  $t$ -step trajectory  $h_t \in \mathcal{H}_{[T]}$ , taking a sub-optimal action  $a \in \mathcal{A} \setminus \{a^*\}$  at step  $t$  in the best-case CMP  $\overline{\mathcal{M}}$  gives a final entropy*

$$H_2^* = \max_{h \in \mathcal{H}_{T-t} \setminus \mathcal{H}_{T-t}^*} H(d_{h_t \oplus h}(\cdot)) \quad \text{s.t. } \mathcal{H}_{T-t}^* = \arg \max_{h \in \mathcal{H}_{T-t}} H(d_{h_t \oplus h}(\cdot))$$

where the maximum is attained by

$$\overline{h}_{T-t} = s_2^* \oplus \left( h_{T-t-1}^* \in \mathcal{H}_{T-t-1}^* \right) \in \arg \max_{h \in \mathcal{H}_{T-t} \setminus \mathcal{H}_{T-t}^*} H(d_{h_t \oplus h}(\cdot)),$$

in which  $s_2^*$  is any state that is the second-closest to a uniform entry in  $d_{h_t \oplus \overline{h}_{T-t}}$ .

*Proof.* The best-case CMP  $\overline{\mathcal{M}}$  is designed such that taking a sub-optimal action  $a \in \mathcal{A} \setminus \{a^*\}$  at step  $t$  minimally decreases the final entropy. Especially, instead of reaching at step  $t + 1$  an optimal state  $s^*$ , i.e., a state that maximally balances the state visits of the final trajectory, the agent is drawn to the second-to-optimal state  $s_2^*$ , from which it gets back on track on the optimal trajectory for the remaining steps. Note that visiting  $s_2^*$  cannot lead to the optimal final entropy, achieved when  $s^*$  is visited at step  $t + 1$ , due to the sub-optimality of action  $a$  at step  $t$  and Assumption 1.  $\square$

**Instantaneous Regret** Although the objective (2) is non-additive across time steps, we can still define a notion of *pseudo-instantaneous regret* by comparing the regret-to-go of two subsequent time steps. In the following, we provide the definition of this expected pseudo-instantaneous regret along with lower and upper bounds to the regret suffered by an optimal Markovian policy.

**Definition B.3** (Expected Pseudo-Instantaneous Regret). *Consider a policy  $\pi \in \Pi$  interacting with a CMP over  $T - t$  steps starting from the trajectory  $h_t$ . We define the expected pseudo-instantaneous regret of  $\pi$  at step  $t$  as  $r_t(\pi) := \max(0, \mathcal{R}_{T-t}(\pi, h_t) - \mathcal{R}_{T-t-1}(\pi, h_{t+1}))$ .*

**Corollary B.4.** *Let  $\pi_M \in \Pi_M$  be a Markovian policy such that  $\pi_M \in \arg \max_{\pi \in \Pi_M} \mathcal{E}(\pi)$  on a CMP  $\mathcal{M}$ . Then, for any  $h_t \in \mathcal{H}_{[T]}$ , it holds  $\underline{r}_t(\pi_M) \leq r_t(\pi_M) \leq \overline{r}_t(\pi_M)$  such that*

$$\begin{aligned} \underline{r}_t(\pi_M) &= \max \left( 0, H^*(\mathcal{V}_t(\pi_M) - \mathcal{V}_{t+1}(\pi_M)) - H_2^* \mathcal{V}_t(\pi_M) + H_* \mathcal{V}_{t+1}(\pi_M) \right), \\ \overline{r}_t(\pi_M) &= \max \left( 0, H^*(\mathcal{V}_t(\pi_M) - \mathcal{V}_{t+1}(\pi_M)) - H_* \mathcal{V}_t(\pi_M) + H_2^* \mathcal{V}_{t+1}(\pi_M) \right), \end{aligned}$$

where

$$\mathcal{V}_t(\pi_M) := \frac{1}{\pi_M(a^*|s_t)} \text{Var}_{h_{s_t} \sim p_t^{\pi_{NM}}} \left[ \mathbb{E} \left[ \mathcal{B}(\pi_{NM}(a^*|h_{s_t})) \right] \right].$$

*Proof.* From the Definition B.3, we have that  $r_t(\pi_M) = \mathcal{R}_{T-t}(\pi_M) - \mathcal{R}_{T-t-1}(\pi_M)$ . Recall that

$$\underline{\mathcal{R}}_{T-t}(\pi) = \mathcal{V}_t(\pi)(H^* - H_2^*), \quad \overline{\mathcal{R}}_{T-t}(\pi) = \mathcal{V}_t(\pi)(H^* - H_*),$$

from Lemma 4.6. Then, we can write

$$\underline{r}_t(\pi_M) \geq \underline{\mathcal{R}}_{T-t}(\pi_M) - \overline{\mathcal{R}}_{T-t-1}(\pi_M) = H^*(\mathcal{V}_t(\pi_M) - \mathcal{V}_{t+1}(\pi_M)) - E_2^* \mathcal{V}_t(\pi_M) + H_* \mathcal{V}_{t+1}(\pi_M),$$

and

$$\overline{r}_t(\pi_M) \leq \overline{\mathcal{R}}_{T-t}(\pi_M) - \underline{\mathcal{R}}_{T-t-1}(\pi_M) = H^*(\mathcal{V}_t(\pi_M) - \mathcal{V}_{t+1}(\pi_M)) - H_* \mathcal{V}_t(\pi_M) + H_2^* \mathcal{V}_{t+1}(\pi_M).$$

□

### B.3. Proofs of Section 5

**Theorem 5.4.**  $\Psi_0$  is NP-hard.

*Proof.* To prove the theorem, it is sufficient to show that there exists a problem  $\Psi_c \in \text{NP-hard}$  so that  $\Psi_c \leq_p \Psi_0$ . We show this by reducing 3SAT, a well-known NP-complete problem, to  $\Psi_0$ . To derive the reduction we consider two intermediate problems, namely  $\Psi_1$  and  $\Psi_2$ . Especially, we aim to show that the following chain of reductions hold:

$$\Psi_0 \geq_m \Psi_1 \geq_p \Psi_2 \geq_p \text{3SAT}$$

First, we define  $\Psi_1$  and we prove that  $\Psi_0 \geq_m \Psi_1$ . Informally,  $\Psi_1$  is the problem of finding a reward-maximizing Markovian policy  $\pi_M \in \Pi_M$  w.r.t. the entropy objective (2) encoded through a reward function in a convenient POMDP  $\widetilde{\mathcal{M}}_\Omega^R$ . We can build  $\widetilde{\mathcal{M}}_\Omega^R$  from the CMP  $\mathcal{M}$  similarly as the extended MDP  $\widetilde{\mathcal{M}}_T^R$  (see Section 2 and the proof of Lemma 4.4 for details), except that the agent only access the observation space  $\widetilde{\Omega}$  instead of the extended state space  $\widetilde{\mathcal{S}}$ . In particular, we define  $\widetilde{\Omega} = \mathcal{S}$  (note that  $\mathcal{S}$  is the state space of the original CMP  $\mathcal{M}$ ),  $\widetilde{O}(\widetilde{o}|\widetilde{s}) = \widetilde{s}_{-1}$ , and the reward function  $\widetilde{R}$  assigns value 0 to all states  $\widetilde{s} \in \widetilde{\mathcal{S}}$  such that  $|\widetilde{s}| \neq T$ , otherwise (if  $|\widetilde{s}| = T$ ) the reward corresponds to the entropy value of the state visitation frequencies induced by the trajectory codified through  $\widetilde{s}$ .

Then, the reduction  $\Psi_0 \geq_m \Psi_1$  works as follows. We denote as  $\mathcal{I}_{\Psi_i}$  the set of possible instances of problem  $\Psi_i$ . We show that  $\Psi_0$  is harder than  $\Psi_1$  by defining the polynomial-time functions  $\psi$  and  $\phi$  such that any instance of  $\Psi_1$  can be rewritten through  $\psi$  as an instance of  $\Psi_0$ , and a solution  $\pi_{\text{NM}}^* \in \Pi_{\text{NM}}$  for  $\Psi_0$  can be converted through  $\phi$  into a solution  $\pi_M^* \in \Pi_M$  for the original instance of  $\Psi_1$ .

$$\begin{array}{ccc} \mathcal{I}_{\Psi_1} & \xrightarrow{\psi} & \mathcal{I}_{\Psi_0} \\ & & \downarrow \\ \pi_M^* & \xleftarrow{\phi} & \pi_{\text{NM}}^* \end{array}$$

The function  $\psi$  sets  $\mathcal{S} = \widetilde{\Omega}$  and derives the transition model of  $\mathcal{M}$  from the one of  $\widetilde{\mathcal{M}}_\Omega^R$ , while  $\phi$  converts the optimal solution of  $\Psi_0$  by computing

$$\pi_M^*(a|o, t) = \sum_{ho \in \mathcal{H}_o} p_t^{\pi_{\text{NM}}^*}(ho) \pi_{\text{NM}}^*(a|ho) \quad (5)$$

where  $\mathcal{H}_o$  stands for the set of histories  $h \in \mathcal{H}_t$  ending in the observation  $o \in \Omega$ . Thus, we have that  $\Psi_0 \geq_m \Psi_1$  holds. We now define  $\Psi_2$  as the policy existence problem w.r.t. the problem statement of  $\Psi_1$ . Hence,  $\Psi_2$  is the problem of determining whether the value of a reward-maximizing Markovian policy  $\pi_M^* \in \arg \max_{\pi \in \Pi_M} \mathcal{J}_{\widetilde{\mathcal{M}}_\Omega^R}(\pi)$  is greater than 0. Since computing an optimal policy in POMDPs is in general harder than the relative policy existence problem (Lusena et al., 2001, Section 3), we have that  $\Psi_1 \geq_p \Psi_2$ .

For the last reduction, i.e.,  $\Psi_2 \geq_p \text{3SAT}$ , we extend the proof of Theorem 4.13 in (Mundhenk et al., 2000), which states that the policy existence problem for POMDPs is NP-complete. In particular, we show that this holds within the restricted class of POMDPs defined in  $\Psi_1$ .

The restrictions on the POMDPs class are the following:



1. The reward function  $R(s) \geq 0$  only in the subset of states reachable in  $T$  steps, otherwise  $R(s) = 0$ ;
2.  $|\tilde{S}| = \tilde{S} = |\tilde{\Omega}|^T$ .

Both limitations can be overcome in the following ways:

1. It suffices to add states with deterministic transitions so that  $T = m \cdot n$  can be defined a priori, where  $T$  is the number of steps needed to reach the state with positive reward through every possible path. Here  $m$  is the number of clauses, and  $n$  is the number of variables in the 3SAT instance, as defined in (Mundhenk et al., 2000);
2. The POMDPs class defined by  $\Psi_1$  is such that  $\tilde{S} = |\tilde{\Omega}|^T$ . Noticing that the set of observations corresponds with the set of variables and that from the previous point  $T = m \cdot n$ , we have that  $|\tilde{\Omega}|^T = n^{m \cdot n}$ , while the POMDPs class used by the proof hereinabove has  $\tilde{S} = m \cdot n^2$ . Notice that  $n \geq 2$  and  $m \geq 1$  implies that  $n^{m \cdot n} \geq m \cdot n^2$ . Moreover, notice that every instance of 3SAT has  $m \geq 1$  and  $n \geq 3$ . Hence, to extend the proof to the POMDPs class defined by  $\Psi_1$  it suffices to add a set of states  $\tilde{S}_p$  s.t.  $R(s) = 0 \forall s \in \tilde{S}_p$ .

Since the chain  $\Psi_0 \geq_m \Psi_1 \geq_p \Psi_2 \geq_p$  3SAT holds, we have that  $\Psi_0 \geq_p$  3SAT. Moreover, since 3SAT  $\in$  NP-complete, we can conclude that  $\Psi_0$  is NP-hard.  $\square$